

Shrinkage Effect in Ancestral Maximum Likelihood

Elchanan Mossel*
 Sebastien Roch[†]
 Mike Steel[‡]

February 8, 2008

Abstract

Ancestral maximum likelihood (AML) is a method that simultaneously reconstructs a phylogenetic tree and ancestral sequences from extant data (sequences at the leaves). The tree and ancestral sequences maximize the probability of observing the given data under a Markov model of sequence evolution, in which branch lengths are also optimized but constrained to take the same value on any edge across all sequence sites. AML differs from the more usual form of maximum likelihood (ML) in phylogenetics because ML averages over all possible ancestral sequences. ML has long been known to be statistically consistent – that is, it converges on the correct tree with probability approaching 1 as the sequence length grows. However, the statistical consistency of AML has not been formally determined, despite informal remarks in a literature that dates back 20 years. In this short note we prove a general result that implies that AML is statistically inconsistent. In particular we show that AML can ‘shrink’ short edges in a tree, resulting in a tree that has no internal resolution as the sequence length grows. Our results apply to any number of taxa.

*Email: mossel@stat.berkeley.edu. Depts. of Statistics and Computer Science, U.C. Berkeley. Supported by an Alfred Sloan fellowship in Mathematics, by NSF grants DMS-0528488, DMS-0548249 (CAREER), and by DOD ONR grant N0014-07-1-05-06.

[†]Email: Sebastien.Roch@microsoft.com. Theory Group, Microsoft Research.

[‡]Email: m.steel@math.canterbury.ac.nz. Biomathematics Research Centre, University of Canterbury, Christchurch, New Zealand

1 Introduction

Markov models of site substitution in DNA are the basis for most methods for inferring phylogenies (evolutionary trees) from aligned sequence data. The usual approach is maximum likelihood (ML) which seeks the tree and branch lengths that maximizes the probability of generating the observed data under a Markov process. In the simplest setting one assumes that sites evolve independently and identically, and that the extant sequences (data) label the leaves of the tree – for background on phylogenetics and ML see [9]. ML is computationally complicated, and even the problem of finding the optimal branch lengths exactly on a fixed tree has unknown complexity. In ML one considers all possible ancestral sequences that could have existed within the tree, and averages each such ‘scenario’ by its probability. An alternative is to simply consider a single choice of ancestral sequences that has the highest probability – this is a variant of ML that was introduced in 1987 by Barry and Hartigan [3] under the name ‘most parsimonious likelihood’, and which later was renamed *ancestral maximum likelihood* (AML) (see e.g. [1]). The computational complexity of AML is slightly easier than ML, in that given the tree and either the optimal branch lengths or the optimal ancestral sequences, the other ‘unknown’ (ancestral sequences or branch length) is readily determined (see eg. [2]). The method can be viewed as being, in some sense, intermediate between ML and a primitive cladistic method, maximum parsimony (MP), which seeks the tree and ancestral sequences that minimizes the total number of sites substitutions required to describe the data. Indeed, AML would select the same trees as MP if one further constrained AML so that each edge had the same branch length, as shown in [10].

The recent interest in AML has sprung from computational complexity considerations. Firstly, AML seemed to provide a promising route by which to show that the problem of reconstructing an ML tree from sequences is NP-hard [1, 6]. It turned out that the NP-hardness of ML can be established directly, without invoking AML [15], however the relative computational simplicity of AML over ML suggests it may provide an alternative strategy for reconstructing large trees.

Nevertheless, it is important to know whether the desirable statistical properties of ML carry over to methods such as AML. In particular ML has long been known to be statistically consistent as a way of estimating tree topologies – that is, as the sequence length grows, the probability that ML will reconstruct the tree that generated the sequences tends to 1. It has also been known (since 1978) that more primitive methods, such as MP, can be statistically inconsistent [8].

However the statistical consistency of AML is unclear, since the standard

Wald-style conditions required to prove consistency (in particular a fixed parameter space that does not grow with the size of the data) does not apply. Thus, one may suspect that AML might be inconsistent, and indeed remarks in the literature have suggested this could be the case (see [4], [11]). However the absence of a sufficient condition to prove consistency does not constitute proof of inconsistency, and the purpose of this short note is to formally show that AML is statistically inconsistent. More precisely we show that AML tends to ‘shrink’ short edges in a tree, and this can result in the collapse of the interior edges (and any short pendant edges) to produce a star tree.

The results in this paper rely on probability arguments, based on expansions of the entropy function, and combinatorial properties of minimal sets of edges that separate each pair of leaves in a tree.

1.1 Problem Statement

CFN model We define $[n] = \{0, \dots, n-1\}$ and we deal with the *Cavender-Farris-Neyman (CFN) model* [5, 7, 13].

Definition 1 (CFN model) We are given a tree $T = (V, E)$ on n leaves labelled $[n]$ and an assignment of edge probabilities $\mathbf{p} : E \rightarrow (0, 1/2)$. A realization of the model is obtained as follows: choose any vertex as a root; pick a state for the root uniformly at random in $\{0, 1\}$; moving away from the root, each edge e flips the state of its ancestor with probability p_e . We denote by X the (random) state at the leaves obtained in this manner. We write $X \sim \text{CFN}(T, \mathbf{p})$.

Ancestral Maximum Likelihood We consider two equivalent formulations of the *Ancestral Maximum Likelihood problem*. The second version is obtained by setting

$$p_e = \frac{d_e}{k}, \tag{1}$$

for all e in the first version [1].

Definition 2 (AML, Version 1) The Ancestral Maximum Likelihood (AML) problem can be stated as follows. Given a set of n binary sequences of length k , find a tree $T = (V, E)$ on n leaves, an assignment $\mathbf{p} : E \rightarrow [0, 1/2]$ of edge probabilities, and an assignment of sequences $\boldsymbol{\lambda} : V \rightarrow \{0, 1\}^k$ to the vertices such that:

1. The sequences at the leaves under $\boldsymbol{\lambda}$ are exactly the sequences from S ;

2. The quantity

$$\mathcal{L}(T, \mathbf{p} \mid \boldsymbol{\lambda}) = -\log_2 \left(\prod_{e \in E} p_e^{d_e} (1 - p_e)^{k - d_e} \right),$$

is minimized, where

$$d_{u,v} = \|\lambda_u - \lambda_v\|_1.$$

Definition 3 (AML, Version 2 [1]) The Ancestral Maximum Likelihood (AML) problem can alternatively be stated as follows. Given a set of n binary sequences of length k , find a tree T on n leaves and an assignment of sequences $\boldsymbol{\lambda} : V \rightarrow \{0, 1\}^k$ to the vertices such that:

1. The sequences at the leaves under $\boldsymbol{\lambda}$ are exactly the sequences from S ;
2. The quantity

$$\mathcal{H}(T \mid \boldsymbol{\lambda}) = \sum_{e \in E} H \left(\frac{d_e}{k} \right),$$

is minimized, where recall that the entropy function is

$$H(p) = -p \log_2 p - (1 - p) \log_2 (1 - p),$$

for $0 \leq p \leq 1$.

Consistency A phylogeny estimator $\Phi = \{(\Phi_n^{(k)})_{n,k \geq 1}\}$ is a collection of mappings from sequences to trees, that is,

$$\Phi_n^{(k)} : \mathcal{B}_n^{(k)} \rightarrow \mathcal{T}_n,$$

where $\mathcal{B}_n^{(k)}$ is the set of all assignments of the form

$$\mathcal{B}_n^{(k)} = \{\boldsymbol{\mu} \mid \boldsymbol{\mu} : [n] \rightarrow \{0, 1\}^k\},$$

and \mathcal{T}_n is the set of all trees on n leaves labelled by $[n]$. Let $\mathbb{X} = \{X_1, X_2, \dots\}$ with $X_j : [n] \rightarrow \{0, 1\}$ for $n \geq 1$. For all $k \geq 1$, we denote by $\boldsymbol{\mu} = \boldsymbol{\mu}_{\mathbb{X}}^{(k)}$ the assignment in $\mathcal{B}_n^{(k)}$ such that $(\mu_v)_j = (X_j)_v$ for all $v \in [n]$ and $j = 1, \dots, k$.

Definition 4 (Consistency) A phylogeny estimator Φ is said to be (statistically) consistent if for all n , all trees $T = (V, E) \in \mathcal{T}_n$, and all edge probability assignments $\mathbf{p} : E \rightarrow (0, 1/2)$, it holds that

$$\Phi_n^{(k)}(\boldsymbol{\mu}_{\mathbb{X}}^{(k)}) \rightarrow T,$$

almost surely as $k \rightarrow +\infty$, where $\mathbb{X} = \{X_1, X_2, \dots\}$ with X_1, X_2, \dots independently generated by $\text{CFN}(T, \mathbf{p})$.

1.2 Main Result

Let Φ_{AML} be the *AML phylogeny estimator* for AML Version 1, where all edges e with $p_e = 0$ have been contracted and all edges e with $p_e = 1/2$ have been removed. (Break ties arbitrarily.)

Theorem 1 (AML Is Not Consistent) *For all $n \geq 1$ and each tree $T = (V, E) \in \mathcal{T}_n$, there is a $\beta > 0$ and a shrinkage zone $\mathcal{Q}_T = \prod_{e \in E} I_e$ such that $|I_e| > \beta$ for all e and if $\mathbf{p} \in \mathcal{Q}_T$, Φ_{AML} returns a star rooted at 0 in the limit $k \rightarrow +\infty$ on the dataset $\mathbb{X} = \{X_1, \dots, X_k\}$ with X_1, \dots, X_k independently generated by $\text{CFN}(T, \mathbf{p})$.*

The phenomenon described in Theorem 1 is illustrated in Fig. 1. We note that our result does not imply the stronger statement that AML is “positively misleading” since we can think of the rooted star as the correct tree T where several edges are set to $p_e = 0$. Note however that the solution is highly degenerate since the star can be obtained in this way from any tree. In other words, in the shrinkage zone, AML provides no information about the internal structure of the tree even with infinitely long sequences.

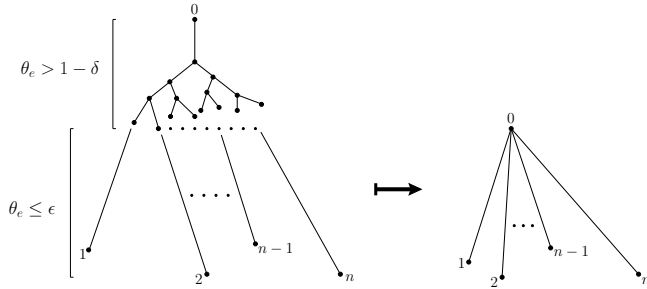


Figure 1: The shrinkage effect: For the tree on the left, AML will reconstruct the star tree (right) from sufficiently long sequences

1.3 Organization

We begin with some preliminary remarks in Section 2. The proof of Theorem 1 can be found in Section 3.

2 Preliminaries

2.1 Solution Properties

Fixed Extension Let $T \in \mathcal{T}_n$. For an assignment of sequences $\boldsymbol{\mu} \in \mathcal{B}_n^{(k)}$ and $1 \leq j \leq k$, we call $\chi : [n] \rightarrow \{0, 1\}$ with $\chi_u = (\mu_u)_j$ for all $u \in [n]$ the j -th character in $\boldsymbol{\mu}$. We write $\chi \in \boldsymbol{\mu}$ if there is j such that χ is the j -th character in $\boldsymbol{\mu}$. We also denote by $\chi^\#$ the number of characters in $\boldsymbol{\mu}$ equal to χ . An extension of a character χ is a mapping $\bar{\chi} : V \rightarrow \{0, 1\}$ such that $\bar{\chi}_v = \chi_v$ for all $v \in [n]$. We denote by $\mathcal{V}(\chi)$ the set of all extensions of χ on T . Let $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$. The mapping then defines an extension for all characters simultaneously by setting $(\bar{\chi}_f)_v = \chi_v$ for all $v \in [n]$ and $(\bar{\chi}_f)_v = f(\chi)_v$ for all $v \in V - [n]$. We show next that AML is in fact equivalent to finding such an f , which can significantly reduce the size of the problem for large k . For a set of n binary sequences $\boldsymbol{\mu} \in \mathcal{B}_n^{(k)}$ and a tree $T = (V, E) \in \mathcal{T}_n$, we denote by $\bar{\boldsymbol{\mu}}_f$ the extension of $\boldsymbol{\mu}$ to V by applying f as above to every character in $\boldsymbol{\mu}$.

Definition 5 (AML, Version 3) Given a set of n binary sequences $\boldsymbol{\mu} \in \mathcal{B}_n^{(k)}$, find a tree $T \in \mathcal{T}_n$ and a mapping $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$ such that the quantity

$$\mathcal{H}(T \mid \bar{\boldsymbol{\mu}}_f) = \sum_{e \in E} H\left(\frac{d_e}{k}\right),$$

is minimized.

Proposition 1 (AML, Version 3) There is always a solution of AML Version 1 and 2 of the form $\boldsymbol{\lambda} = \bar{\boldsymbol{\mu}}_f$ for some $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$.

Proof: Note that

$$\begin{aligned} \mathcal{L}(T, \mathbf{p} \mid \boldsymbol{\lambda}) &= -\log_2 \left(\prod_{e \in E} p_e^{d_e} (1 - p_e)^{k - d_e} \right), \\ &= -k \sum_{e \in E} \log_2(1 - p_e) - \\ &\quad \sum_{j=1}^k \sum_{(u,v) \in E} \mathbb{1}\{(\lambda_u)_j \neq (\lambda_v)_j\} \log_2 \frac{p_e}{1 - p_e}. \end{aligned}$$

For fixed \mathbf{p} , since \mathcal{L} “decomposes” in j , it is always possible to take the same extension for each character appearing in $\boldsymbol{\mu}$ without affecting optimality. Then, we can choose the optimal \mathbf{p} as in [1] to obtain the result. ■

Limit Problem Let $T = (V, E) \in \mathcal{T}_n$. Assume as in Theorem 1 that we are given a dataset $\mathbb{X} = \{X_1, X_2, \dots\}$ with X_1, X_2, \dots i.i.d. $\text{CFN}(T, \mathbf{p})$. Fix $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$. Let $X \sim \text{CFN}(T, \mathbf{p})$ and denote by $Y = \bar{X}_f$ the extension of X under f . Also, let $\bar{\mu}_{\mathbb{X},f}^{(k)}$ be the extension of $\mu_{\mathbb{X}}^{(k)}$ under f . By the Law of Large Numbers, as $k \rightarrow +\infty$, the quantity $\mathcal{H}(T \mid \bar{\mu}_{\mathbb{X},f}^{(k)})$ converges almost surely to

$$\mathbb{H}_{X,T}(f) = \sum_{e \in E} H(Y_e),$$

where, for $e = (u, v)$, Y_e is the indicator that $Y_u \neq Y_v$, and $H(Y_e)$ is the entropy of Y_e , that is,

$$H(Y_e) = H(\mathbb{P}[Y_u \neq Y_v]).$$

Note that, by Proposition 1, even as $k \rightarrow +\infty$ there are only a constant number of mappings f to consider. We say that f is $\mathbb{H}_{X,T}$ -optimal if f minimizes $\mathbb{H}_{X,T}(f)$ over all $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$. The minimum need not be unique.

Definition 6 (Expected AML) *Given a random variable X taking values in $\{0, 1\}^{[n]}$, find a tree $T = (V, E) \in \mathcal{T}_n$ and a mapping $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$ such that the quantity*

$$\mathbb{H}_{X,T}(f) = \sum_{e \in E} H(Y_e),$$

is minimized, where $Y = \bar{X}_f$.

By the previous remarks and (1), to prove Theorem 1 it suffices to show:

Theorem 2 (Optimal Assignment) *Let $T' = (V', E') \in \mathcal{T}_n$ and let $X \sim \text{CFN}(T', \mathbf{p})$. Then there is a $\beta > 0$ and a shrinkage zone $\mathcal{Q}_T = \prod_{e \in E} I_e$ such that $|I_e| > \beta$ for all e and for all $T = (V, E) \in \mathcal{T}_n$, the unique $\mathbb{H}_{X,T}$ -optimal $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$ assigns to all internal nodes of V the value at leaf 0 under all characters, that is,*

$$f(x) = (x_0, \dots, x_0),$$

for all $x \in \{0, 1\}^{[n]}$.

2.2 Minimal Isolating Sets

Definition In preparation for our proof of Theorem 2, we will need the following notion which is studied in [12].

Definition 7 (Isolating Set) Let $T = (V, E)$ be a tree. A subset S of E is called an isolating set for T if for any two leaves u, v there exists an edge $e \in S$ on the path connecting u and v .

The following result is proved in [12].

Proposition 2 (Minimal Isolating Set) The size of a minimal isolating set on an n -leaf tree is $n - 1$.

We will also need:

Proposition 3 (One Leaf Per Component) Let T be a tree on n leaves and let S be a minimal isolating set on T . Consider the forest F obtained from T by removing all edges in S . Then, each component of F contains exactly one leaf of T .

Proof: If a component of F contains two leaves, then these cannot be isolated under S , a contradiction. On the other hand, if a component T' of F does not contain a leaf, then every edge adjacent to T' in T is in fact in S . But then one can remove one of these edges without losing the isolating property of S , contradicting the minimality of S . ■

Minimally Isolating f Let $T = (V, E) \in \mathcal{T}_n$ and $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$. We denote by $S_f \subseteq E$ the set of edges $e = (u, v)$ such that there is $x \in \{0, 1\}^{[n]}$ with $f(x)_u \neq f(x)_v$.

Definition 8 (Minimally Isolating f) We say that f is minimally isolating for T if S_f is a minimal isolating set of T .

2.3 Random cluster parameterization

We will sometimes require a different ('random cluster') parameterization of the CFN model. Let $T \in \mathcal{T}_n$ and $\mathbf{p} \in [0, 1]^E$. (Note that we allow p_e in $[0, 1]$.) We let

$$\theta_e = 1 - 2p_e,$$

for all $e \in E$. The main property we will use is the following well-known identity. For two leaves u, v in T , let $\text{Path}_T(u, v)$ be the set of edges on the path between u and v .

Proposition 4 (Path Probability) *Let $T = (V, E) \in \mathcal{T}_n$ and $\mathbf{p} \in [0, 1]^E$. Assume $X \sim \text{CFN}(T, \mathbf{p})$. Let u, v be two leaves of T . Then we have*

$$\mathbb{P}[X_u \neq X_v] = \frac{1}{2} \left(1 - \prod_{e \in \text{Path}_T(u, v)} \theta_e \right).$$

3 Proof

In this section, we prove Theorem 2 from which Theorem 1 follows. The proof has two components:

1. [Reduction to Minimal Isolating Sets] We first show that for any random variable $X \in \{0, 1\}^{[n]}$ close enough to uniform and any tree $T \in \mathcal{T}_n$, the $\mathbb{H}_{X, T}$ -optimal f 's are minimally isolating for T .
2. [Rooted Star is Optimal] Second, we show that if X above is $\text{CFN}(T', \mathbf{p})$ for some $T' \in \mathcal{T}_n$ with $p_e \approx 1/2$ if e is adjacent to $\{1, \dots, n-1\}$ and $p_e \approx 0$ otherwise, then for all $T \in \mathcal{T}_n$ the unique $\mathbb{H}_{X, T}$ -optimal f assigns the value at 0 to all internal nodes.

Throughout, $n \geq 1$ is fixed.

3.1 Reduction to Minimal Isolating Sets

We prove the following:

Proposition 5 (Reduction to Minimal Isolating Sets) *There exists $\varepsilon > 0$ (depending on n) such that the following hold. Let X be any random variable taking values in $\{0, 1\}^{[n]}$ with $H(X) \geq n - \varepsilon$ and let T be any tree in \mathcal{T}_n . If f is $\mathbb{H}_{X, T}$ -optimal, then f is minimally isolating for T .*

Proof: We make a series of claims.

Claim 1 (Reduction to Uniform) *For all $\delta > 0$ there exists $\varepsilon = \varepsilon(\delta) > 0$ such that if X is a $\{0, 1\}^{[n]}$ -random variable with*

$$H(X) \geq n - \varepsilon,$$

and $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$ then

$$|\mathbb{H}_{X, T}(f) - \mathbb{H}_{U, T}(f)| \leq \delta, \tag{2}$$

where U is the uniform distribution on $\{0, 1\}^{[n]}$. Therefore, it suffices to prove Proposition 5 for those f that are $\mathbb{H}_{U,T}$ -optimal.

Proof: The entropy of $\{0, 1\}^{[n]}$ -random variables is maximized uniquely at $H(U) = n$. The first part of the result follows by continuity of $H(X)$ and $\mathbb{H}_{X,T}(f)$ in the distribution of X .

For the second part, take $\delta > 0$ small enough such that for all f, f' , we have

$$\mathbb{H}_{U,T}(f) > \mathbb{H}_{U,T}(f') \implies \mathbb{H}_{U,T}(f) > \mathbb{H}_{U,T}(f') + 2\delta. \quad (3)$$

(Recall that there are only finitely many f 's for fixed n .) Take $\varepsilon > 0$ such that the first part holds. Then it follows that if f is $\mathbb{H}_{X,T}$ -optimal then it must be $\mathbb{H}_{U,T}$ -optimal. We argue by contradiction. Assume there are f, f' such that $\mathbb{H}_{X,T}(f) \leq \mathbb{H}_{X,T}(f')$ but $\mathbb{H}_{U,T}(f) > \mathbb{H}_{U,T}(f')$. By (3), we have

$$\mathbb{H}_{U,T}(f) > \mathbb{H}_{U,T}(f') + 2\delta, \quad (4)$$

which implies $\mathbb{H}_{X,T}(f) > \mathbb{H}_{X,T}(f')$ by (2), a contradiction. ■

Claim 2 (Minimizer) *If f is $\mathbb{H}_{U,T}$ -optimal then $\mathbb{H}_{U,T}(f) = n - 1$. Moreover, denoting $Y = \bar{U}_f$ we have that $\{Y_0, (Y_e)_{e \in E}\}$ are mutually independent.*

Proof:

Upper Bound We first show that there is f such that $\mathbb{H}_{U,T}(f) \leq n - 1$. Let S be a minimal isolating set for T . Define f by letting $f(x)_u = f(x)_v$ for all edges (u, v) not in S . By Proposition 3, this uniquely defines f . Letting $Y = \bar{U}_f$ it is immediate to check that

$$\mathbb{H}_{U,T}(Y) = \sum_{e \in E} H(Y_e) = \sum_{e \in S} H(Y_e) \leq n - 1,$$

by Proposition 2.

Lower Bound For any $f : \{0, 1\}^{[n]} \rightarrow \{0, 1\}^{V-[n]}$ with $Y = \bar{U}_f$, we have

$$\begin{aligned} n = H(U) &= H(\{(Y_v)_{v \in [n]}\}) = H(\{Y_0, (Y_e)_{e \in E}\}) \\ &\leq H(Y_0) + \sum_{e \in E} H(Y_e) \leq 1 + \sum_{e \in E} H(Y_e), \end{aligned}$$

where we have used that $\{(Y_v)_{v \in [n]}\}$ and $\{Y_0, (Y_e)_{e \in E}\}$ are deterministic functions of each other. Furthermore, the first inequality holds to equality if and only if $\{Y_0, (Y_e)_{e \in E}\}$ are mutually independent. ■

We are ready to conclude the proof of Proposition 5. Let f be $\mathbb{H}_{U,T}$ -optimal with $Y = \bar{U}_f$. Let u, v be any two leaves of T . We have by the previous claim that $(Y_e)_{e \in \text{Path}_T(u,v)}$ are mutually independent. Since Y_u and Y_v are independent uniform $\{0, 1\}$ it must be that there is an edge $e \in \text{Path}_T(u, v)$ with $H(Y_e) = 1$. Indeed, define $p_e = \mathbb{P}[Y_e = 1]$ and $\theta_e = 1 - 2p_e$. Then by Proposition 4 we have

$$0 = 1 - 2\mathbb{P}[Y_u \neq Y_v] = \prod_{e \in \text{Path}_T(u,v)} \theta_e,$$

which implies that at least one $\theta_e = 0$. Let S' be the set of all edges where $H(Y_e) = 1$. Then we have shown that S' is an isolating set. Note furthermore that if $e \in S_f$ then $H(Y_e) \geq H(2^{-n}) > 0$. From f 's optimality we obtain

$$n - 1 = \mathbb{H}_{U,T}(f) \geq |S'| + |S_f \setminus S'|H(2^{-n}).$$

Therefore we must have $S_f = S'$ and $|S'| = n - 1$ which implies that S_f is a minimal isolating set as needed. ■

3.2 The Rooted Star is Optimal

Let $T = (V, E) \in \mathcal{T}_n$ and S a minimal isolating set of T . Let T^0 be the tree obtained from T by contracting all edges not in S . By Proposition 3, T^0 is an n -node tree where each node (leaf or internal) is (uniquely) labelled by a leaf of T . Let \mathcal{T}_n^0 be all such trees on n nodes. By Proposition 5, the AML phylogeny estimator is among \mathcal{T}_n^0 . Note that for $T \in \mathcal{T}_n^0$ the only possible extension is trivially $f = \mathbb{1}$ since there are no unlabelled internal vertices.

Proposition 6 (Rooted Star is Optimal) *Let $T = (V, E) \in \mathcal{T}_n$. Let W be the set of leaf edges of T , except the edge pendant at 0. Then for $\varepsilon, \delta > 0$ sufficiently small the following holds. Assume $X \sim \text{CFN}(T, \mathbf{p})$ with corresponding random cluster parameterization satisfying $0 < \theta_e \leq \varepsilon$ for all $e \in W$ and $1 > \theta_e > 1 - \delta$ for all $e \notin W$. Then, among all trees $T' \in \mathcal{T}_n^0$, the star rooted at 0 uniquely minimizes $\mathbb{H}_{X,T'}(\mathbb{1})$ for all δ sufficiently small.*

Proof: We assume that δ and ε are small enough so that they satisfy

$$(n - 1)(1 - \delta)^{2n-4} > n - 2,$$

and

$$\varepsilon^2 < (n-1)(1-\delta)^{2n-4} - (n-2). \quad (5)$$

Let $T' = (V', E') \in \mathcal{T}_n^0$ and $f = \mathbb{1}$ with corresponding variables $(Y_0, \{Y_e\}_{e \in E'})$ where $Y_0 = X_0$ and $Y_{u,v} = \mathbb{1}\{X_u \neq X_v\}$. Let $e = (u, v)$ be an edge in T' . In particular, note that u and v are leaves of T . Let $p_{u,v}$ be the probability that u and v disagree and let $\theta_{u,v} = 1 - 2p_{u,v}$. We will use the following Taylor expansion of the entropy around $1/2$

$$H\left(\frac{1-\tau}{2}\right) = 1 - \left(\frac{\log_2 e}{2}\right) \tau^2 + O(\tau^4).$$

Note further that

$$H(Y_e) = H(p_{u,v}) = H\left(\frac{1-\theta_{u,v}}{2}\right).$$

As ε approaches 0, $p_{u,v}$ goes to $1/2$. Therefore, by Proposition 4, up to smaller order terms we want to find $T' = (V', E')$ in \mathcal{T}_n^0 that maximizes

$$\Theta(T') := \sum_{e'=(u,v) \in E'} \prod_{e \in \text{Path}_T(u,v)} \theta_e^2.$$

If T' has an edge e' between two leaves neither of which is 0, then e' makes a contribution of at most ε^4 to $\Theta(T')$ since $\text{Path}_T(u, v)$ crosses two edges in W . Therefore, by (5),

$$\begin{aligned} \Theta(T') &\leq (n-2)\varepsilon^2 + \varepsilon^4 \\ &< (n-1)(1-\delta)^{2n-4}\varepsilon^2, \end{aligned}$$

where we have used that T' has exactly $n-1$ edges and each edge $e'' \in E'$ makes a contribution of at most ε^2 since $\text{Path}_T(u, v)$ contains at least one edge in W . On the other hand, the star rooted at 0, which we denote by T^* , is the only tree in \mathcal{T}_n^0 which does not include an edge between two leaves neither of which is 0. In that case, we get

$$\Theta(T^*) \geq (n-1)(1-\delta)^{2(n-2)}\varepsilon^2,$$

where we have used that any path between 0 and another leaf in T contains at most $n-2$ edges not in W (since $|E| \leq 2n-3$ and $|W| = n-1$) and exactly one edge in W . Taking ε small enough gives the result. ■

4 Concluding remarks

It would be interesting to extend our results beyond the 2-state case. We note in particular that for the symmetric r -state model, with $r > 2$, the equivalent formulation of the AML problem given in Definition 3 does not apply. Indeed, it is easy to check that, instead, one needs to minimize

$$\mathcal{H}'(T \mid \lambda) = \sum_{e \in E} H\left(\frac{d_e}{k}\right) + \log_2(r-1) \sum_{e \in E} \frac{d_e}{k}.$$

The second term on the r.h.s.—a parsimony “correction”—may lead to a different behavior when $r > 2$.

We thank Peter Ralph for sharing his recent, independent results [14] regarding the structure of the optimal solution in the 2-state case (similarly to [2]) as well as a number of simulations on 4-taxon trees.

References

- [1] L. Addario-Berry, B. Chor, M. Hallett, J. Lagergren, A. Panconesi, T. Wareham, “Ancestral Maximum Likelihood of Evolutionary Trees is Hard,” *Jour. of Bioinformatics and Comp. Biology*, vol. 2, no. 2, pp. 257-271, 2004.
- [2] N. Alon, B. Chor, F. Pardi, A. Rapoport, “Approximate maximum parsimony and ancestral maximum likelihood,” to appear in *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 2008.
- [3] D. Barry, J. Hartigan, “Statistical analysis of hominoid molecular evolution,” *Stat. Sci.*, 2, pp. 191-207, 1987.
- [4] D. Barry, J. Hartigan, “Rejoinder [on Statistical analysis of hominoid molecular evolution],” *Stat. Sci.*, vol. 2, pp. 209-210, 1987.
- [5] J. A. Cavender, “Taxonomy with confidence,” *Math. Biosci.*, vol. 40, no. 3-4, 1978.
- [6] B. Chor, T. Tuller, “Finding the maximum likelihood tree is hard,” in *Proc. 9th Annual International Symposium on Research in Computational Biology (RECOMB 2005)*, 2005.

- [7] J. S. Farris, "A probability model for inferring evolutionary trees", *Syst. Zool.*, vol. 22, no. 4, pp. 250-256, 1973.
- [8] J. Felsenstein, 1978. "Cases in which parsimony or compatiability methods will be positively misleading", *Syst. Biol.*, 27, pp. 401–410.
- [9] J. Felsenstein, *Inferring phylogenies*, Sinauer, Sunderland, MA, 2004.
- [10] N. Goldman, 1990. "Maximum likelihood inference of phylogenetic trees, with special reference to a Poisson process of DNA substitution and to parsimony analysis", *Syst. Zool.* 39, pp. 345–361.
- [11] P. A. Goloboff, "Parsimony, likelihood, and simplicity," *Cladistics*, vol. 19, pp. 91-103, 2003.
- [12] V. Moulton, M. Steel, "Peeling phylogenetic 'oranges'," *Adv. in Appl. Math.*, vol. 33, no. 4, pp. 710-727, 2004.
- [13] J. Neyman, "Molecular studies of evolution: a source of novel statistical problems," in *Statistical desicion theory and related topics*, (eds. S. S. Gupta and J. Yackel), pp. 1-27, Academic Press, New York, 1971.
- [14] P. Ralph, in preparation, 2008.
- [15] S. Roch, "A Short Proof that Phylogenetic Tree Reconstruction by Maximum Likelihood is Hard," *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, vol. 3, no. 1, pp. 92-94, 2006.